

Word Length Frequency and Distribution in English: Observations, Theory, and Implications for the Construction of Verse Lines

Hideaki AOYAMA and John CONSTABLE

*Faculty of Integrated Human Studies,
Kyoto University, Kyoto 606-8501, Japan*

`aoyama@phys.h.kyoto-u.ac.jp`

`john@ic.h.kyoto-u.ac.jp`

January, 1998

Abstract

Recent observations in the theory of verse and empirical metrics have suggested that constructing a verse line involves a pattern-matching search through a source text, and that the number of found elements (complete words totaling a specified number of syllables) is given by dividing the total number of words by the mean number of syllables per word in the source text. This paper makes this latter point explicit mathematically, and in the course of this demonstration shows that the word length frequency totals in English output are distributed geometrically (previous researchers reported an adjusted Poisson distribution), and that the sequential distribution is random at the global level, with significant non-randomness in the fine structure. Data from a corpus of just under two million words, and a syllable-count lexicon of 71,000 word-forms is reported. The pattern-matching theory is shown to be internally coherent, and it is observed that some of the analytic techniques described here form a satisfactory test for regular (isometric) lineation in a text.

Keywords: word length, metrics

1 Introduction

The making of metrical verse lines is a pattern-matching exercise, and in this respect contrasts sharply with the generation of ordinary language output. The composer must first generate source text, and then search through it for elements which form, or may be accumulated to form, the desired pattern. This pattern is usually specified in relation to a selection from the available surface features of the language, those in English verse, for example, being based principally on the patterning of stressed and unstressed syllables to form arrangements of beats and offbeats (Attridge 1982, Attridge 1995). In earlier work Constable (Constable 1997) has observed that these rules lead to the implication of a simpler and more readily studied rule, of which a composer is usually unaware, namely that a line should consist of complete words totaling n syllables, where n can be a range. This may be put as a rule thus:

‘In every consecutive section of n syllables there must be only complete words’.

The composition of word strings to fit a line definition rule of this type involves finding sequences of complete words totaling n syllables, or constructing them from sequences of less than n syllables. In either case the activity is a pattern-matching search for a target, and the frequency of the target in the source language is of crucial interest to the composer, since a smaller number of found elements will more severely restrict the communicative options open. That is to say, the more targets are present in the source text, the more probable it is that the composer will find pieces that function adequately with regard to his or her purpose. Therefore, in order to facilitate line composition, authors are expected to take whatever action they can to increase the frequency of target elements. An effective way of doing this is to reduce the mean number of syllables per word in the source language being searched. Recent work in empirical metrics (Constable 1997: 182) has demonstrated the relation between mean word length and target frequency empirically from small samples of English prose, and has claimed that the average frequency of a target object (a sequence of words totaling n syllables) is given by the total number of words divided by the mean number of syllables per word (and this leads to the prediction that when composing in verse authors will tend to choose shorter words than they would otherwise have selected, a point also explored in (Constable 1997)). The mathematical reformulation required to substantiate and explain the factors underlying this effect involves detailed observations regarding the frequency and sequential word length distributions typical of English, and we will now turn to this task.

2 Global Structure

In approaching word-length data, it is of crucial importance to form an analytic strategy suitable for the purpose in hand. Given a distribution, either of frequencies or probabili-

ties, it is possible to fit it with any degree of accuracy to a function, as long as one has a large number of functions, each with a large enough number of parameters. Previous research on word length frequency undertaken by the Göttingen Word Length Project (see Best and Altmann 1996 for an overview of the project and a bibliography; prominent papers in the project include: Becker 1996, Dittrich 1996, Frischen 1996, Riedemann 1996, Röttger 1996, Wimmer *et. al.* 1996, Wimmer 1994, Ziegler 1996, Zuse 1996) has collected data relating to word length frequencies, usually from rather small samples, and then used software, the Altmann Fitter, to compute a best fit description, which in most cases proves to be an adjusted Poisson distribution (Best and Altmann 1996). Such a fitting procedure is guaranteed to work: the number of words with more than seven syllables is small, and they are infrequent in output, so there are, normally, only approximately seven data points, and, therefore, it will not be difficult to find an accurate fit if we have several hundred functions, each with several parameters, from which to choose. Mathematically, if there are seven data items, a general function with 7 parameters would suffice for a *perfect* fit.

We stress that this is not relevant: Although such a fit may serve certain technical purposes, it does not bring any insight into or understanding of the nature of the language under consideration. It is only when we can find a fit with a much smaller number of parameters than that of the data, that we can move towards abstraction and understanding. Any deviation from this simple ansatz should be regarded as a subtle variation from the basic finding.

Methodology of this type is quite common in exact science. For example, in high-energy physics, where the purpose of experiments, say proton-electron collisions, is not to fit the experimental data with one of several thousand functions, but to grasp the nature of protons and their constituent quarks. Once we know the fundamental properties, further exploration of any deviation in the experimental data opens the way to advance our understanding.

Our research is directed with this principle in mind, and we have therefore chosen not to fit the distributions using a large set of functions and parameters. Rather, we will first extract a simple property that describes the overall, global, structure of the data, from which we aim to obtain a deeper knowledge of the English language, if only in its syllabic organization. With this in hand we can then proceed to the study of deviations from this global structure, that is to the fine structure of the data.

2.1 The Corpus

We have analyzed the word-length structure of almost two million words of prose by various authors as listed in Table 1. The texts were chosen with no other view than convenience, Constable having marked them in the course of other research. We acknowledge that a more principled choice would be desirable, and indeed hope to undertake such work ourselves, but we believe that this corpus is sufficient for present purposes. The syllabic

data were obtained by using a simple marking program, constructed by Constable, which reads text and uses a custom-built lexicon to determine the syllabic count of each word form. When new word forms are encountered, the program requests human intervention, and the form is added to the lexicon. Consistency of syllable counting was ensured by the fact that only one user (Constable) has been responsible for building the lexicon, which at the time of writing contains 71,666 items. Abbreviations were expanded, and numbers were counted as if they were pronounced as hyphenated words (1,920 = one-thousand-nine-hundred-and-twenty), with the exception of years and dates which were counted in their normal pronounced form (1920 = nineteen twenty). None of these categories are, as it happens, frequent in our corpus.

2.2 Frequency and Probabilities

In order to present the method of analysis in definite terms, let us introduce several mathematical notations and facilities. We denote the syllabic data obtained as described above by a series of integers N_i ($i = 1, 2, 3, \dots, I$), where N_i is the number of syllables in the i -th word and I the total number of words in the data.

The number of sequences in a series N_i which satisfy a line definition rule of the type introduced in the previous section is represented as follows;

$$n = \sum_{i=\ell}^{m \pmod{I}} N_i. \quad (1)$$

The upper limit of this sum implies that the line definition rule is applied with a ‘periodic boundary condition’; namely, the data is treated as a circle by connecting the end of data with the beginning. This is a technical definition justifiable by its utility in the following mathematical treatment. Alternatively, one could use a Dirichlet boundary condition, in which one simply terminates the data sequence at $i = I$. These boundary conditions, however, do not significantly affect the results as long as the data size is large, which is true for all the data we have analyzed.

It is straightforward to count the number $L_{n,k}$ which match Eq.(1) for $k = m - \ell$ words: The numbers $L_{n,1}$ are obtained simply by counting the numbers equal to n among the series (N_1, N_2, \dots, N_I) . Next, the numbers $L_{n,2}$ are obtained by counting similarly for $(N_1 + N_2, N_2 + N_3, \dots, N_I + N_1)$, $L_{n,3}$ from $(N_1 + N_2 + N_3, N_2 + N_3 + N_4, \dots, N_I + N_1 + N_2)$, and so forth. By definition, the following identity is satisfied:

$$\sum_{n=1}^{\infty} L_{n,k} = I. \quad (2)$$

Since there are no zero-syllable words in English,

$$L_{n,k} = 0 \quad \text{if} \quad n < k. \quad (3)$$

The quantity we are interested in is the number of sequences matching the line definition rule for *any* number of words, which is given by the following:

$$L_n = \sum_{k=1}^n L_{n,k}. \quad (4)$$

This counting algorithm has been coded in *Mathematica* by Aoyama and has been found to work much faster than the original algorithm used by Constable (Constable 1997:181). In Table 2 we give the partial list of $L_{n,k}$ and L_n obtained for all the data listed in Table 1.

For theoretical reasons it is best to deal with quantities independent of the data size. Therefore, we introduce the following normalized quantities:

$$P_{n,k} \equiv \frac{L_{n,k}}{I}, \quad Q_n \equiv \frac{L_n}{I}. \quad (5)$$

Due to the identity (2), the following is satisfied:

$$\sum_{n=1}^{\infty} P_{n,k} = 1. \quad (6)$$

In this sense, a set of $P_{n,k}$ of a given k defines a probability distribution. On the other hand, Q_n does not have this property. We call Q_n a (normalized) ‘frequency’. Corresponding to Eq.(4), we have the following relation:

$$Q_n = \sum_{k=1}^n P_{n,k}. \quad (7)$$

In Fig.1 we give the plot of $P_{n,k}$ and Q_n for all the data listed in Table 1. As is seen in Fig.1, the most remarkable global feature of the frequency distribution Q_n is its flatness, that is, its independence from n . In order to analyze this structure, we may define an idealized constant distribution $\bar{Q}_n = q$, where $q = 0.720316$ is the average value of the actual distribution Q_n for $n = 1 \sim 30$. In the following section we will examine what lies behind this constant distribution.

2.3 Random-Ordering Hypothesis

As a working hypothesis we might assume that the word-length series is randomly ordered, or more accurately, that

‘the number of syllables in a word is independent of the number of syllables in preceding words.’

In other words, this hypothesis suggests that there is no correlation between the syllable-count values in the data series. This random-ordering hypothesis allows us to express $P_{n,k}$ in terms of $P_{n,1}$, the probability of a word having n syllables (hereafter we denote this

quantity by $p_n = P_{n,1}$). For example, two consecutive words that satisfy the two-syllable line definition rule can be obtained by having two one-syllable words in a row. The number of one-syllable words is Ip_1 , and according to the random-ordering hypothesis the probability of having a one-syllable word after a one-syllable word is not affected by the first word having one-syllable, and therefore is p_1 . Thus the number of 2-syllable lines of this form is given by $Ip_1 \times p_1$. Dividing this by the total number of words, I , we obtain the normalized frequency $P_{2,2}$:

$$P_{2,2} = p_1^2. \quad (8)$$

For larger n and k , combinatoric considerations must be addressed. For example, a three-syllable line can be created by having a two-syllable word and a one-syllable word in sequence, or vice versa. Counting all possibilities, we obtain,

$$P_{3,2} = 2p_1p_2. \quad (9)$$

Some of the other relations are listed in Table 3. The general expression for $P_{n,k}$ can be obtained in a straightforward manner, but is complicated in written form, and can be handled most simply by the use of generating functions.

We define a generating function $P_k(x)$ to represent $P_{n,k}$ for $n = 1 \sim \infty$ by the following:

$$P_k(x) \equiv \sum_{n=1}^{\infty} P_{n,k} x^n. \quad (10)$$

Knowledge of all $P_{n,k}$ is equivalent to knowing the behaviour of $P_k(x)$ near the origin $x = 0$, as $P_{n,k}$ can be expressed as the n -th order derivative of $P_k(x)$ at $x = 0$:

$$P_{n,k} = \frac{1}{n!} \frac{d^n P_k}{dx^n}(0). \quad (11)$$

The normalization condition (6) of $P_{n,k}$ is expressed as $P_k(1) = 1$. The various moments of n (expectation values of powers of n) can be expressed in terms of derivatives of $P_k(x)$ at $x = 1$. For example, the average $\langle n \rangle$ and the standard deviation σ of n are given by the following:

$$\langle n \rangle_k \equiv \sum_{n=1}^{\infty} n P_{n,k} = P'_k(1), \quad (12)$$

$$\begin{aligned} \sigma_k &\equiv \sqrt{\langle (n - \langle n \rangle_k)^2 \rangle_k} = \sqrt{\langle n^2 \rangle_k - \langle n \rangle_k^2} \\ &= \sqrt{P''_k(1) + P'_k(1) - (P'_k(1))^2}. \end{aligned} \quad (13)$$

We also define a generating function $Q(x)$ as follows:

$$Q(x) \equiv \sum_{n=1}^{\infty} Q_n x^n. \quad (14)$$

In terms of these generating functions, the relation (7) is written as

$$Q(x) = \sum_{k=1}^{\infty} P_k(x). \quad (15)$$

A relation similar to Eq.(11) holds also for Q_n .

The general expression of $P_{n,k}$ in terms of p_n induced by the random ordering condition can be summarized very simply. In terms of the generating functions it is expressed as follows:

$$P_k(x) = P_1(x)^k. \quad (16)$$

We note that the normalization is trivial in the above equation: $P_k(1) = P_1(1)^k = 1$. The relation (16) leads to the following expression of the generating function $Q(x)$:

$$Q(x) = \sum_{k=1}^{\infty} P_k(x) = \sum_{k=1}^{\infty} P_1(x)^k = \frac{P_1(x)}{1 - P_1(x)}. \quad (17)$$

Thus the reason for introducing the random-ordering hypothesis becomes evident: If that hypothesis is valid, the features of the frequency distribution Q_n can be explained by the features of the one-word probability distribution p_n by using the relation (17).

We will now turn to the verification of the random-ordering hypothesis. In Table 4 we list the number of m -syllable words following immediately after n -syllable words. The corresponding probability distribution is plotted in Fig. 2. From this figure, we readily observe that this distribution is almost independent from the value of m . Therefore we confirm that the random-ordering hypothesis is valid to a reasonable degree of accuracy.¹⁾

2.4 Single Word Probability

Now that the random-ordering hypothesis is confirmed, we can obtain the probability \bar{p}_n that induces the constant frequency distribution $\bar{Q}_n = q$. The generating function for \bar{Q}_n is as follows:

$$\bar{Q}(x) = \sum_{n=1}^{\infty} qx^n = \frac{qx}{1-x}. \quad (18)$$

By solving Eq. (17) in terms of $P_1(x)$, we obtain,

$$\bar{P}_1(x) = \frac{\bar{Q}(x)}{1 + \bar{Q}(x)} = \frac{qx}{1 - (1-q)x} = \sum_{n=1}^{\infty} q(1-q)^{n-1}x^n. \quad (19)$$

Therefore,

$$\bar{p}_n = q(1-q)^{n-1}, \quad (20)$$

which is the geometric probability distribution. In Fig. 3 we compare the actual probability distribution p_n (dots) and the geometric distribution \bar{p}_n in Eq. (20) (dash-dotted line). As is seen in this plot, the agreement is close. The geometric distribution (20) yields the average number of the syllables per word (mean word length) as follows:

$$\langle n \rangle = \sum_{n=1}^{\infty} n\bar{p}_n = \bar{P}'_1(1) = \frac{1}{q}. \quad (21)$$

This relation was observed earlier by Constable (Constable 1997:182). We stress that our new finding of the relation between the constant distribution \bar{Q}_n and the geometric distribution \bar{p}_n , which we reached through the application of the random-ordering hypothesis gives a sound theoretical basis to this observation.

2.5 Interpretation: Random Segmentation

The two global properties we have found above, the random-ordering and the geometric distribution (20), allow a definite characterization of the word-length data, since these are the properties typical of a system with a given probability of termination at any point: namely, if one assumes that sequences of syllables are constructed such that after any syllable, the end of a word happens with probability q , the above geometric distribution (20) is obtained. Putting this in a slightly different manner, if one has a large number of syllables and word boundaries (spaces) with $(1 - q)$ to q ratio and randomly places them in sequence, the same distribution is obtained. We call this *random segmentation*.

The fact that this geometric distribution is not found in the lexicon itself, which is plotted in Fig. 4, has been noted by other researchers (Wimmer *et. al.* 1996), and provoked explanation in terms of attractors and control cycles in the composition process. Strictly speaking it is beyond the scope of our paper to engage deeply with this question. However, since we believe that these researchers have been misled by the presumption of order in the output distribution, there is some point in observing that hypotheses based on randomness could, in principle, account for the relations between these very different distributions, and explain the stability of the geometric distribution in output.

For example, we might hypothesize that the concept of ‘word’ or ‘word boundary’ is a relatively late (though perhaps prehistoric) analytic category, and has been arrived at by segmenting the verbal output stream in such a way that word boundaries are placed with a fixed probability in relation to syllable boundaries. The resulting word-forms are used to compile a lexicon. Since the sound system of a language is not infinitely extendible, there will be more unique and acceptable disyllabic forms than monosyllables, more trisyllables than monosyllables, and so on. Thus, although in its early stages the lexicon would, obviously, follow the geometric distribution of the output it was drawn from (Fig.3), eventually it would, temporarily, adopt the sort of curve seen in Fig.4.

This is not to suggest that the segmentation of English is fundamentally random, or that ‘words’ have a low linguistic reality, interesting though both speculations are. In line with the data examined so far, our hypothesis merely notes that whatever principle of regular order may be operating elsewhere, perhaps in relation to stress or phonemes, word boundaries and syllable boundaries are related with a fixed probability.

3 Fine Structures

Readers will have noticed the differences between the global structure and the actual distributions. The most notable is the small dip of Q_2 below the average value q seen in Fig.1, and the small differences between $p_{n,m}$ of different m in Fig.2.

When we discuss these differences, we need first to guard against statistical errors. In other words, we first need to see whether these differences are meaningful quantities or can be attributed to statistical fluctuations. Only when the former is more likely, do

we need to study the fine structures that explain these differences. We stress that this discussion of statistical errors is of the first importance. As we see below small data sets, such as those employed by the Göttingen group, do in fact suffer from large statistical errors. Detailed study of such data is either irrelevant or misleading. In the following, we first discuss the handling of statistical errors and then proceed to the discussion of features of the fine structure.

3.1 Statistical Errors

The standard estimate for statistical errors may be applied to the individual probabilities that we deal with this paper. The $3\text{-}\sigma$ error range for the probability $P_{n,k}$ would be,

$$P_{n,k} - 3\sqrt{\frac{P_{n,k}(1 - P_{n,k})}{I}} \sim P_{n,k} + 3\sqrt{\frac{P_{n,k}(1 - P_{n,k})}{I}}. \quad (22)$$

In other words, the true value lies in this range with about 99.7% probability.

The estimate of the error range for the frequencies require more extensive discussion. Such an estimate is made possible by relation (7): It is trivial for Q_1 , as it is actually a probability, $Q_1 = P_{1,1}$. Therefore its standard deviation is given by $\sigma = \sqrt{Q_1(1 - Q_1)}/\sqrt{I}$. The next is $Q_2 = P_{2,1} + P_{2,2}$. The probability $P_{2,1}$ is given by dividing the observed number $L_{2,1}$ of 2-syllable words by the total number of words I . From Fig.3 we see that $P_{2,1} \sim 0.2$, therefore the standard deviation of $L_{2,1}$ can be approximated as $\sigma_{2,1} \sim \sqrt{L_{2,1}}$. Similarly, we approximate that $\sigma_{2,2} \sim \sqrt{L_{2,2}}$. Thus the total standard deviation for the observed number of strings L_2 is given as follows:

$$\begin{aligned} \sigma_2^2 &= \langle L_2^2 \rangle - \langle L_2 \rangle^2 \\ &= L_{2,1} + L_{2,2} = L_2, \end{aligned} \quad (23)$$

where we used the statistical independence of $L_{2,1}$ and $L_{2,2}$. That is, the standard deviation of Q_2 is given simply by $\sqrt{Q_2/I}$, just as if it is a probability by itself. The same is true for the rest of Q_n s.

The estimate of the statistical errors for the average value q of Q_n is simplified because of the flatness of the distribution Q_n . Given the fact that the value of Q_n is independent from n , the statistical errors follow from statistical errors of any *one* value of Q_n , which is almost independent from n , as explained above. Therefore, we simply estimate the standard deviation of q to be of the same order as a typical value of that of Q_n s;²⁾ Namely, we estimate the 3σ -range of q to be between $q \pm 3\sqrt{q/I}$.

3.2 Deviations from the Flat Q_n Distribution

In order to study the small deviation of the frequency Q_n from the flat distribution, we can plot the difference between the actual frequency and the flat distribution itself, $\delta Q_n = Q_n - q$, with vertical bars showing the 3σ error ranges, as in Fig.5.

From this figure, we find that the 2 syllable depression plot is statistically significant, as are other deviations at $n = 1, 3, 4$. These fine deviations can be explained from underlying deviations; namely, (1) deviation from the geometric distribution, (2) deviation from random-ordering. In order to examine these Fig.6 plots (a) $p_n - \bar{p}_n$ (solid line) and (b) $p_{n,1} - \bar{p}_n$ (dotted line). The line (a) is a measure of the deviation from the geometric distribution, while the difference between (a) and (b) is a measure of the deviation from random-ordering. In these figures we find that (1) monosyllabic words have a slightly higher probability than that predicted by the geometric distribution, while disyllabic words have a slightly lower probability, and that (2) in the sequential distribution there is a slightly enhanced probability of a polysyllable after a monosyllable (relative to random-ordering). These are the important, indeed the only, exceptions to the overall randomness in the distributions. These deviations, being small compared to unity, can be mathematically treated as perturbations in a randomly-ordered geometric distribution. In that manner, it is straightforward to show that these small deviations for smaller n do not affect Q_n for large n , thus explaining the fact that deviation of Q_n from the flat distribution is localized to small n .

We have not examined the underlying linguistic causes of these deviations, and are not in a position to do more than speculate. The corpus is predominantly of high status literary writing, mostly nineteenth-century, and of that a substantial portion comes from one author, Henry James. It might therefore be suggested that these deviations are characteristic of an output type, or a period, or even of an author. However, some of the deviations observed overall are consistently found across authors and works, though in the case of Kipling we found a significant depression at $n = 1$ instead of an enhancement. We predict, therefore, that some of these deviations, the depression at $n = 2$ for example, are universal characteristics, while some others will prove to be particular to an author, a work, a genre type, or a period. With regard to the deviation from random sequencing, we suggest that the relation between commonly occurring function terms, which are predominantly monosyllabic, and content terms, which are somewhat more likely to be polysyllabic, is the likeliest explanation.

Whatever the best causal account, it should be recognized that deviations such as these are subtle variations from a strong fundamental trend, and it is not safe to conclude that they are evidence which ‘confirm[s] the assumption of a non-accidental distribution of word lengths’ (Ziegler 1996:73).

3.3 Individual Authors

The global structures we noted in the preceding section, the flatness of the frequency Q_n and the randomness of sequencing are also true for individual authors and works. In other words, it is not a result of averaging over large fluctuations among authors. In Fig.7 and Fig.8 we give the plot of the frequency Q_n and the probability $p_{n,m}$ for George Eliot’s novel *Middlemarch*. The global features are readily apparent from these figures.

The average frequency is $q = 0.69844$ for the Eliot data and is $q = 0.720316$ for all the data in Table 1, and readers may wonder whether this difference is meaningful. As explained in a previous subsection, the standard deviation σ_q of the average frequency q is estimated to be \sqrt{q}/I . This means that $\sigma_q \simeq \sqrt{0.720316}/1977676 \simeq 0.00060$ for the corpus listed in Table 1. Thus at the 3σ -confidence level, the true value of q lies between 0.7185 and 0.7221. Similarly, the 3σ -confidence range of q for the George Eliot data is $0.6940 \sim 0.7029$. Since these ranges do not overlap, we conclude that indeed the mean number of syllables per word in Eliot is significantly larger than that of the corpus. In Fig.9 we give a plot of the values of q and the 3σ ranges of all the authors and the whole corpus.

3.4 Test for Lineation

Apart from the $n = 2$ deviation observable in the prose texts in our corpus, we are aware of one large class of texts which routinely exhibit significant deviations in the Q_n distribution, namely isometrically lineated verse texts. This is hardly surprising. Texts composed in regular lines are by definition ordered with respect to lineation, and this order will be detected by such a procedure as ours. If a poem is composed in lines of ten syllables, for example, then $n = 10$, and all multiples of ten, will be substantially above the flat distribution, and if it is composed in two core line lengths, as limericks and Spenserian stanzas are, then it will exhibit two series of peaks. Since we intend to discuss this matter at greater length elsewhere we will present only one example, the final version of William Wordsworth's *Prelude*, in Fig.10 and Fig.11. This poem, which was completed in 1839 but not published until 1850, is composed in blank verse, that is unrhymed five-beat lines in duple rhythm, with a range of between 9 and 12 syllables per line (in duple rhythm the offbeat position is usually a single syllable, but is sometimes filled with two syllables, or even left unfilled). The poem contains 7,849 lines and 57,570 words, with a mean number of syllables per word of 1.4.

Approximately 77.5 % of the text is composed in ten syllable lines, with 19.4 % being of eleven syllables and 2.3 % of twelve syllables, together with a scattering of other lengths. This degree of concentration into the core line length is by no means abnormal, and if anything it is somewhat more distributed than other texts examined. The Q_n distribution reveals, apart from the expected depression at $n = 2$, significant peaks at ten and all multiples of ten, though the subsequent peaks are of course of lesser size, since long, uninterrupted, runs of ten-syllable lines are rare.

It should be noted that this test is of theoretical rather than practical value. It is unlikely that we will often wish to test in order to detect lineation, since the fact is usually visually evident, or clear from other features such as rhythm. However, as a contribution to the theoretical definition of verse lines and texts, particularly as distinguished from prose, the procedure is of considerable interest. Although it has been long obvious that lineation is not merely 'a visual or typographical fact' but a 'fact of the language'

(Wimsatt and Beardsley 1959:591), to use one well-known formula, there has been, to our knowledge, no conclusive empirical demonstration of the presence of this fact, or any explanation of its character. The deviation from the flat distribution performs both these functions.

4 Conclusion and Comments

Previous research on word length distribution (Becker 1996, Best and Altmann 1996, Dittrich 1996, Frischen 1996, Riedemann 1996, Röttger 1996, Wimmer *et. al.* 1996, Wimmer 1994, Ziegler 1996, Zuse 1996) has attempted to infer significance from the non-geometric curve found, and held that it supports the belief that ‘language is [...] a self-regulating system, which is controlled by the needs of the language community’ (Zuse 1996), or is an organism of interrelated control cycles (Wimmer *et. al.* 1996). Furthermore, these researchers have incautiously borrowed terms from chaos theory, and been, in our view, misled by them. For example, in Wimmer *et. al.* 1996:98 it was claimed on the basis of a handful of data, that ‘the sequence of words is clearly chaotic’, and that the distribution of word length in a text could be explained by reference to ‘attractors’. However, in the current context, there is in fact no chaos, in its mathematical sense, and what we observe in our study is randomness: when the sample size is small, any distribution, height or weight in a human population, or, to mention something fundamentally random, quantum theoretical events, will exhibit large fluctuations, and as the data size grows, the distributions become smoother. We do not rule out the discovery of the sort of order sought by the synergetic linguists, but observe that our findings give little support to its existence in relation to word length. Thus, in approaching frequency data of this type we find ourselves generally in sympathy with those such as Mandelbrot (Mandelbrot 1961) and Li (Li 1992), who are of course studying different linguistic features, in their advocacy of interpretations grounded in randomness, and we are less drawn to positions such as those proposed by Zipf (Zipf 1965:48), where statistical regularities are seen to arise from some deep principle of order.

In conclusion, however, we should like to emphasize that the theory and data outlined here are of more than negative value, or purely self-sufficing interest. Our investigation was derived from empirical observations and hypotheses offered in an earlier paper (Constable 1997) with regard to the construction of verse lines, and those remarks are confirmed by our results. The relation between the mean number of syllables per word and the number of sequences of words totalling a given number of syllables (Constable 1997:182) is dependent on the geometric frequency of word length totals, and the random distribution of the words in the text sequence, which we have shown here to be solid findings. Thus, the apparently arcane facts of word length distribution in English output can be seen to deepen our understanding of one, and a historically very important, area of language output, isometrical verse.

Footnotes

1. We note that these features, flatness of the frequency Q_n and random-ordering are also true for individual authors and works, and are not a result of averaging over them. These issues will be addressed in subsequent sections.
2. One might think that having a number of data points for Q_n would reduce the standard deviation for q by a factor equal to the square root of the number of Q_n . However, since the flat distribution has $\bar{Q}_n = q$ for $n = 1 \sim \infty$, this argument would yield a zero standard deviation, which is clearly wrong.

References

- Attridge, D. (1982). *The Rhythms of English Poetry*. Longman, London.
- Attridge, D. (1995). *Poetic Rhythm: An introduction*. Cambridge University Press, Cambridge.
- Becker, C. (1996). Word Lengths in the Letters of the Chilean Author Gabriela Mistral. *Journal of Quantitative Linguistics* 3.128-131.
- Best, K.-H. and Altmann, G. (1996). Project Report. *Journal Quantitative Linguistics* 3.85-88.
- Constable, J. (1997). Verse Form; A Pilot Study in the Epidemiology of Representations. *Human Nature* 8.171-203.
- Dittrich, H. (1996). Word Length Frequency in the Letters of G. E. Lessing. *Journal of Quantitative Linguistics* 3.26-264.
- Frischen, J. (1996). Word Length Analysis of Jane Austen's Letters. *Journal of Quantitative Linguistics* 3.80-84.
- Li, W. (1992). Random Texts Exhibit Zipf's-Law-Like Word Frequency Distribution. *IEEE Transactions on Information Theory* 38.1842-1845.
- Mandelbrot, B. (1961). On the Theory of Word Frequencies and on Related Markovian Models of Discourse. In *Structure of Language and its Mathematical Aspects*. Proceedings of Symposia in Applied Mathematics. 12.190-219. American Mathematical Society.
- Riedemann, H. (1996). Word Length Distribution in English Press Texts. *Journal of Quantitative Linguistics* 3.265-271.
- Röttger, W. (1996). Distribution of Word Length in Ciceronian Letters. *Journal of Quantitative Linguistics* 3.68-72.
- Wimmer, G., Köhler, R., Grotjahn, R., and Altmann, G. (1996). Towards a Theory of Word Length Distribution. *Journal of Quantitative Linguistics*, 1.98.
- Wimmer, G. and Altmann, G. (1994). The theory of word length: some results and generalizations. *Glottometrika*, 15.
- Wimsatt, W. K., and Beardsley, M. C. (1959). The Concept of Meter: An exercise in abstraction. *PMLA* 74.585-598.
- Ziegler, A. (1996). Word Length Distribution in Brazilian-Portuguese Texts. *Journal of Quantitative Linguistics* 3.73-79.

Zipf, G. K. (1965). *The Psychobiology of Language*. M.I.T. Press, Cambridge Mass.

Zuse, M. (1996). Distribution of Word Length in Early Modern English: Letters of Sir Philip Sidney. *Journal of Quantitative Linguistics* 3.272-6.

<i>Author</i>	<i>Title(s)</i>	<i>Words</i>
Bunyan, John	<i>Pilgrim's Progress</i>	52,504
Eliot, George	<i>Middlemarch</i>	317,827
Frankau, Gilbert	<i>Woman of the Horizon</i> (section; first 67 pages 10 chapters)	24,597
Goldsmith, Oliver	<i>Vicar of Wakefield</i>	63,096
James, Henry	<i>The Altar of the Dead; The Ambassadors; The American; The Aspern Papers; Confidence; Daisy Miller; Death of the Lion; The Europeans; The Figure in the Carpet; The Golden Bowl; An International Episode; Portrait of a Lady; Roderick Hudson; Sacred Fount; Turn of the Screw; Watch and Ward; Washington Square</i>	1,285,041
Kipling, Rudyard	<i>Rewards and Fairies; The Jungle Book</i>	115,602
Milton, John	<i>History of Britain; Colasterion; Martin Bucer</i>	119,009
<i>Total</i>		1,977,676

Table 1: Content of the corpus; authors, titles of the sources and number of words from each author

n	$L_{n,1}$	$L_{n,2}$	$L_{n,3}$	$L_{n,4}$	$L_{n,5}$	L_n
1	1,433,426	0	0	0	0	1,433,426
2	371,500	1,025,719	0	0	0	1,397,219
3	122,179	558,679	733,202	0	0	1,414,060
4	40,314	246,132	611,737	531,686	0	1,429,869
5	9,048	99,647	348,154	583,554	387,684	1,428,087
6	1,082	33,891	169,801	411,842	524,656	1,425,780
7	119	9,790	73,374	238,242	439,613	1,426,115
8	6	2,983	27,502	121,820	293,091	1,426,874
9	2	660	9,832	54,843	171,197	1,426,456
10	0	146	2,902	22,952	88,979	1,426,660
11	0	26	878	8,329	42,419	1,426,218
12	0	3	219	3,012	18,275	1,426,323
13	0	0	60	983	7,464	1,426,066
14	0	0	13	299	2,856	1,425,963
15	0	0	2	91	941	1,425,162
16	0	0	0	16	333	1,425,536
17	0	0	0	4	108	1,425,480
18	0	0	0	2	44	1,424,226
19	0	0	0	0	8	1,424,392
20	0	0	0	1	6	1,425,044
21	0	0	0	0	1	1,425,327
22	0	0	0	0	0	1,425,568
23	0	0	0	0	1	1,425,068
24	0	0	0	0	0	1,424,803
25	0	0	0	0	0	1,424,248
26	0	0	0	0	0	1,424,738
27	0	0	0	0	0	1,424,738
28	0	0	0	0	0	1,424,730
29	0	0	0	0	0	1,424,089
30	0	0	0	0	0	1,424,313

Table 2: Number of strings $L_{n,k}$ and L_n that satisfy the n -syllable line definition rule for $n = 1 \sim 30$ and $k = 1 \sim 5$. The values of $L_{n,k}$ for $k = 6 \sim 30$ are omitted due to space limitations. These figures cover all two-million words of data listed in Table 1.

	$n = 1$	2	3	4	5
$P_{n,1}$	p_1	p_2	p_3	p_4	p_5
$P_{n,2}$	0	p_1^2	$2p_1p_2$	$2p_1p_3 + p_2^2$	$2(p_1p_4 + p_2p_3)$
$P_{n,3}$	0	0	p_1^3	$3p_1^2p_2$	$3(p_1^2p_3 + p_1p_2^2)$
$P_{n,4}$	0	0	0	p_1^4	$4p_1^3p_2 + 6p_1^2p_2^2$
$P_{n,5}$	0	0	0	0	p_1^5

Table 3: Some of the consequences of the Random-Ordering Hypothesis. The probability $P_{n,k}$ is listed at the (k, n) position.

m	$n = 1$	2	3	4	5	6	7	8	9
all	1,433,426	371,500	122,179	40,314	9,048	1,082	119	6	2
1	1,025,719	279,915	90,473	29,731	6,680	811	91	5	1
2	278,764	63,357	20,960	6,733	1,485	183	17	0	1
3	92,302	19,542	7,263	2,422	589	52	8	1	0
4	29,414	6,815	2,710	1,128	217	27	3	0	0
5	6,400	1,617	692	265	67	7	0	0	0
6	745	225	71	30	9	2	0	0	0
7	75	28	10	5	1	0	0	0	0
8	5	1	0	0	0	0	0	0	0
9	2	0	0	0	0	0	0	0	0

Table 4: List of the number of occurrences of n -syllable words after m -syllable words

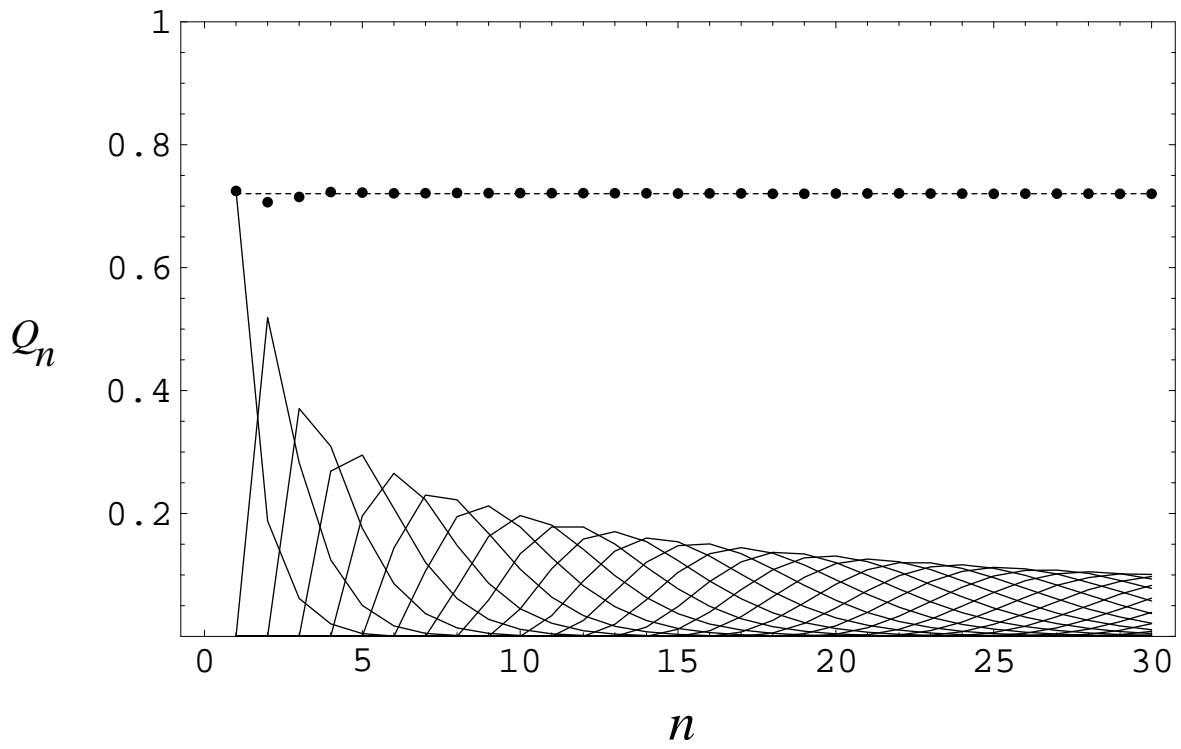


Figure 1: Plot of the normalized frequency Q_n and the probability distributions $P_{n,k}$ for the data in Table 1. The horizontal dashed line shows the average value $q = 0.720316$ for Q_n .

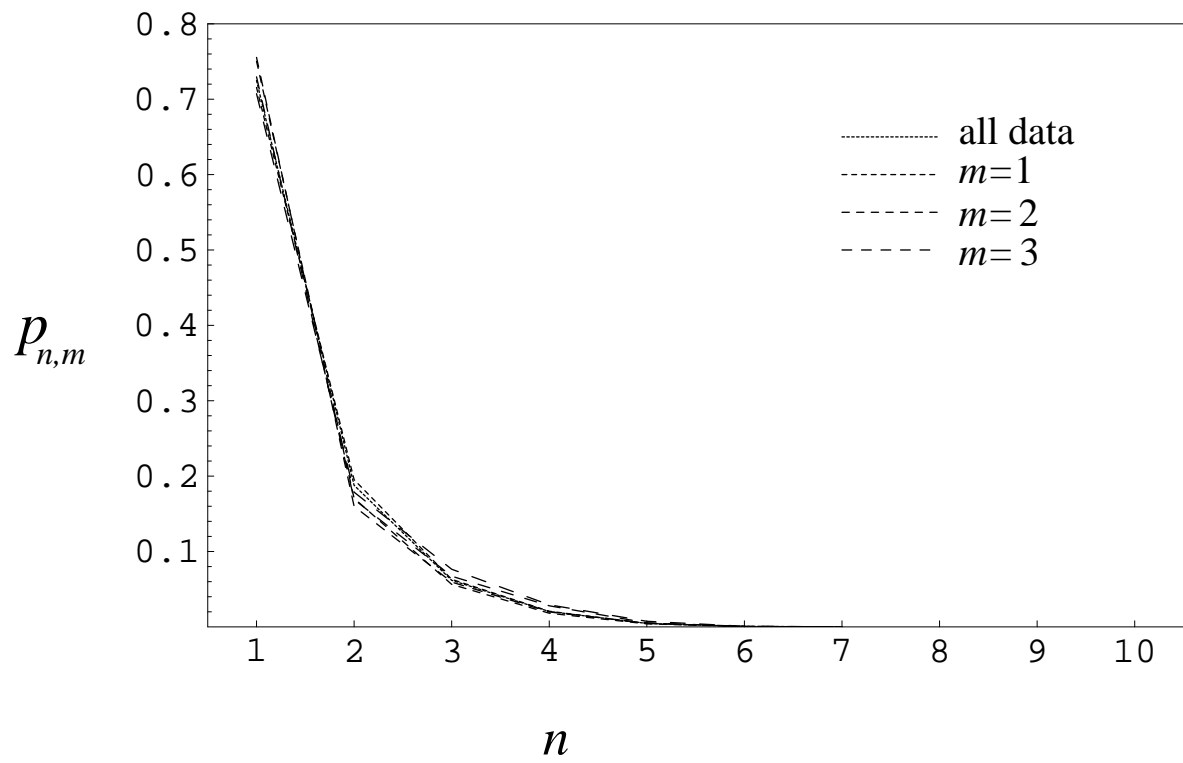


Figure 2: Plot of the probability distribution $p_{n,m}$ for the data in Table 4. The solid line shows p_n , while other lines show $p_{n,m}$.

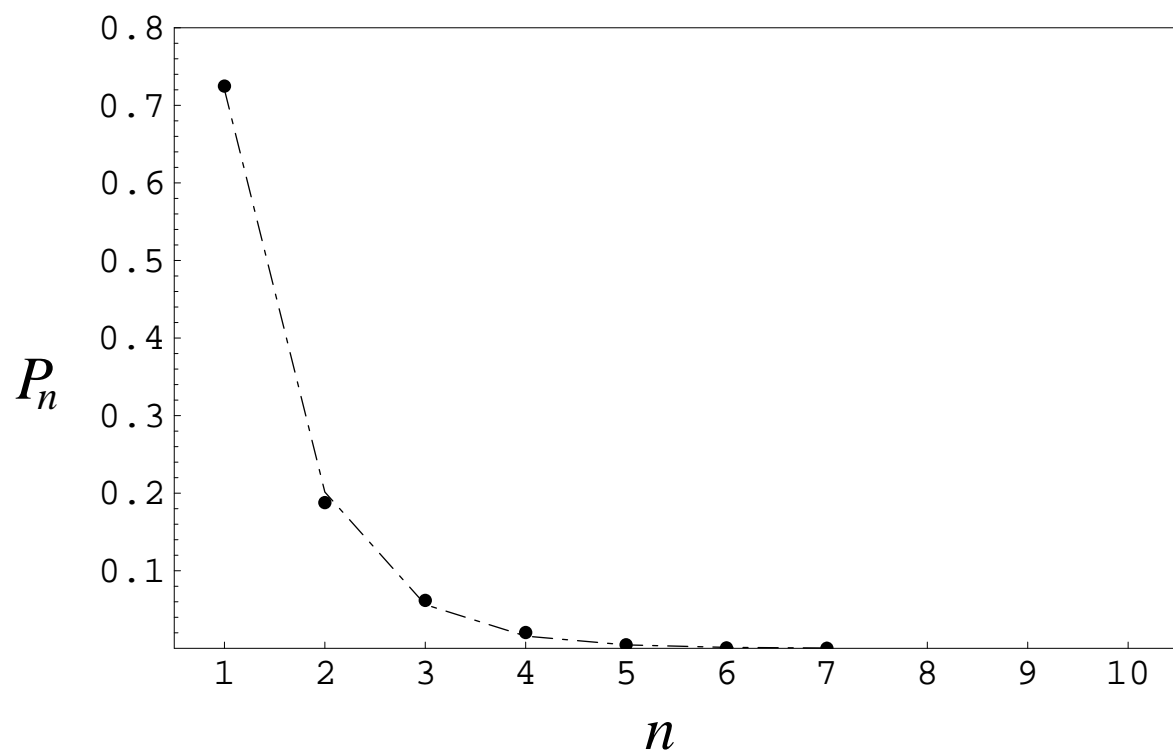


Figure 3: The probability distribution p_n (denoted by dots) and the theoretical prediction \bar{p}_n (20) (denoted by the dash-dotted line)

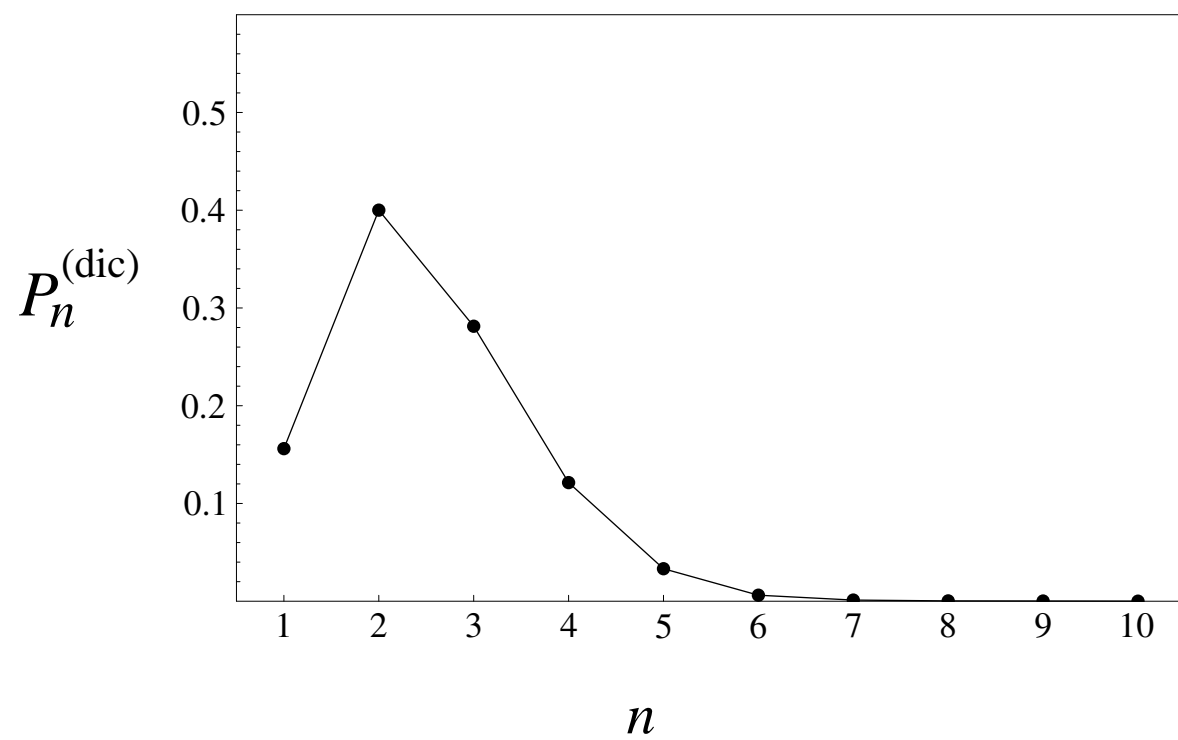


Figure 4: Probability distribution for words with n -syllables in Constable's lexicon

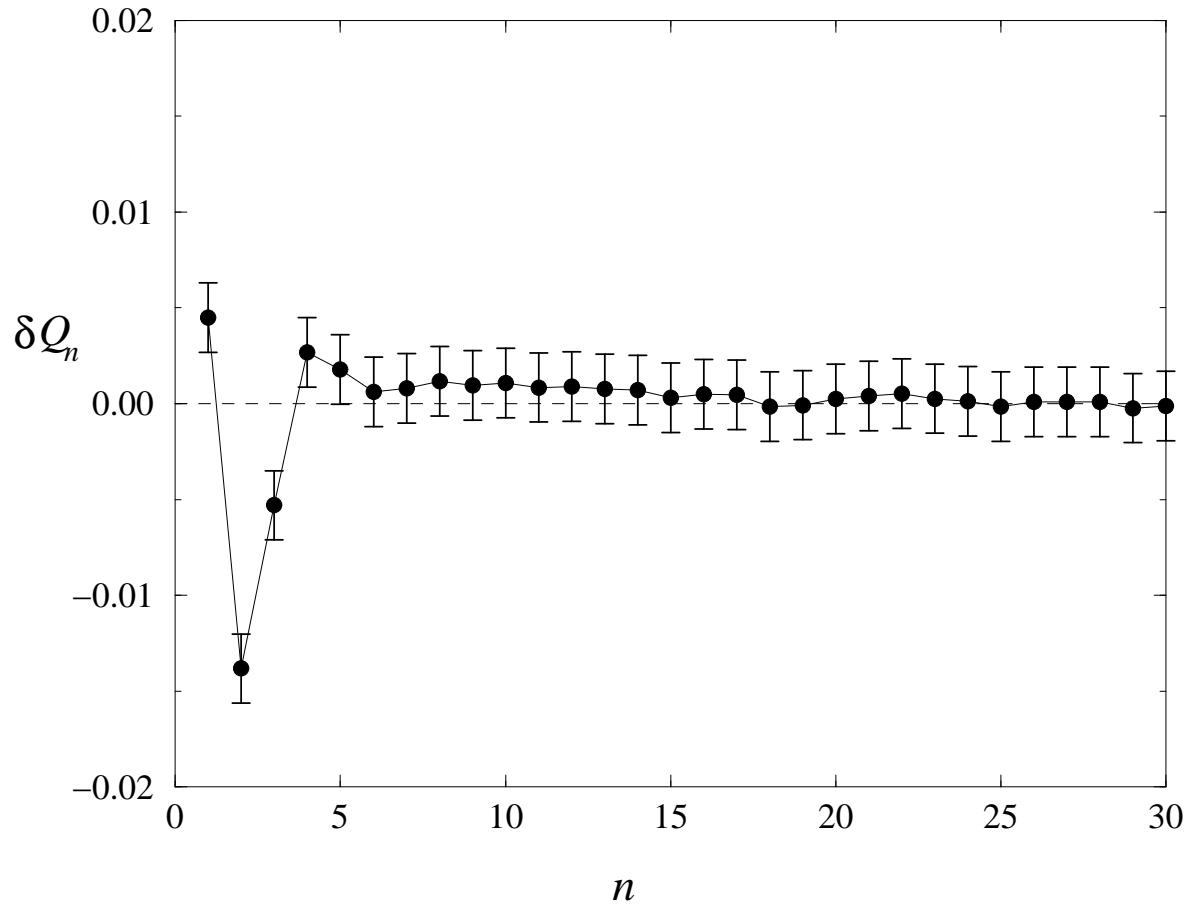


Figure 5: Detail of $\delta Q_n = Q_n - q$. The vertical bars show the 3σ -confidence ranges of statistical errors for each data point. Note that the horizontal range covers only 1/25 of the range of Fig.1.

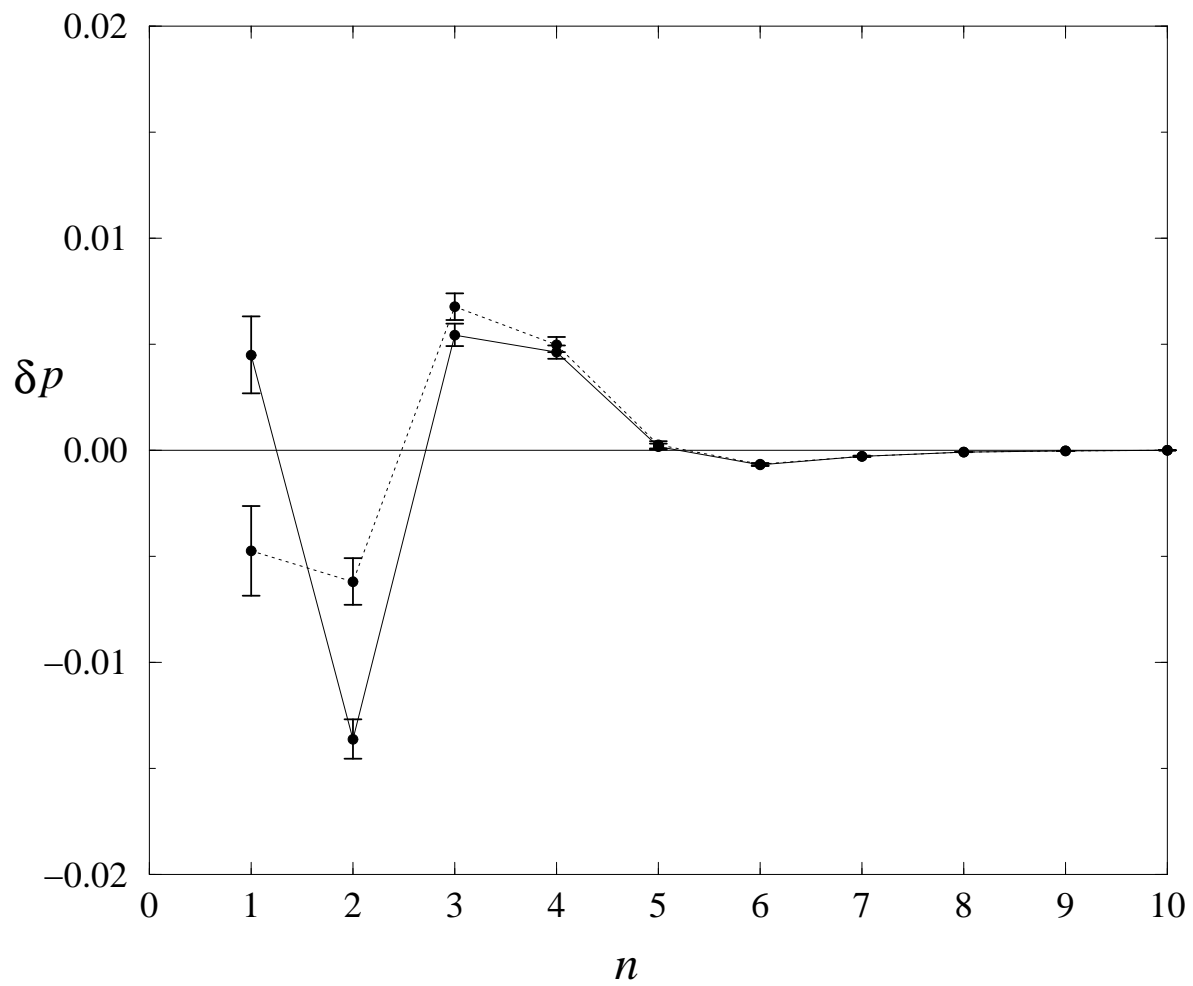


Figure 6: Detail of (a) $p_n - \bar{p}_n$ (solid line) and (b) $p_{n,1} - \bar{p}_n$ (dotted line)

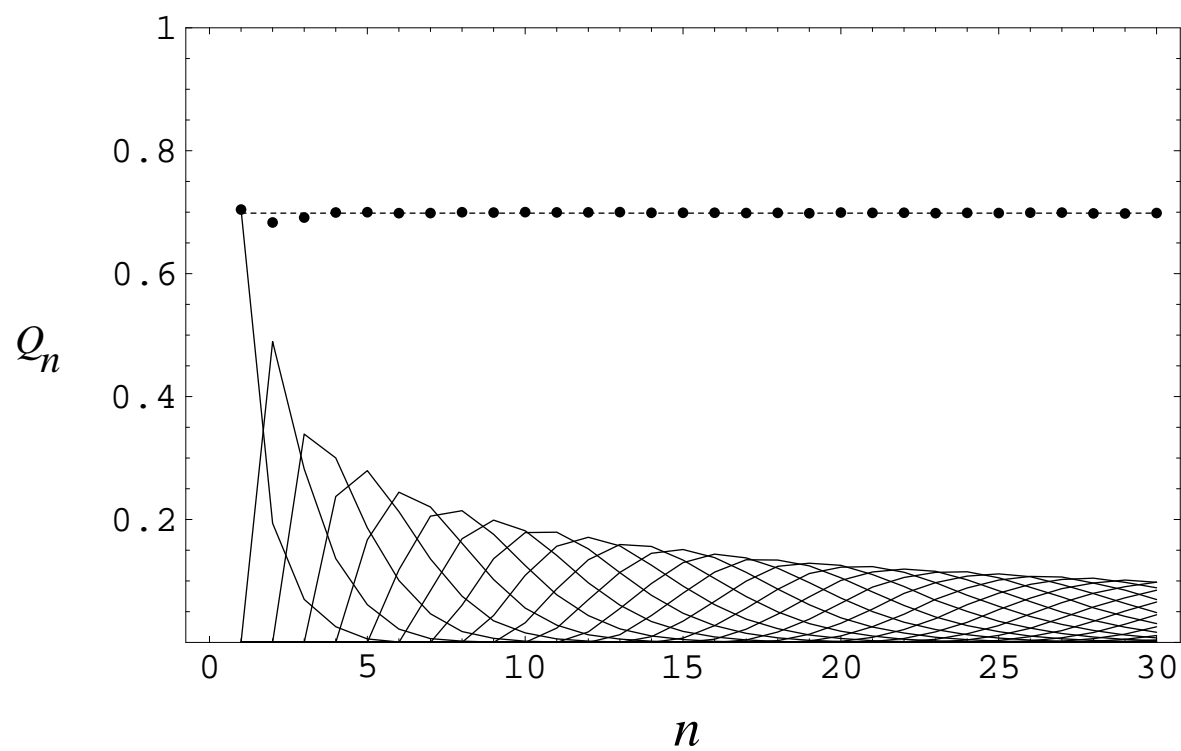


Figure 7: Plot of the normalized frequency Q_n and the probability distributions $P_{n,k}$ for George Eliot, *Middlemarch*. The horizontal dashed line shows the average value $q = 0.69844$ for Q_n .

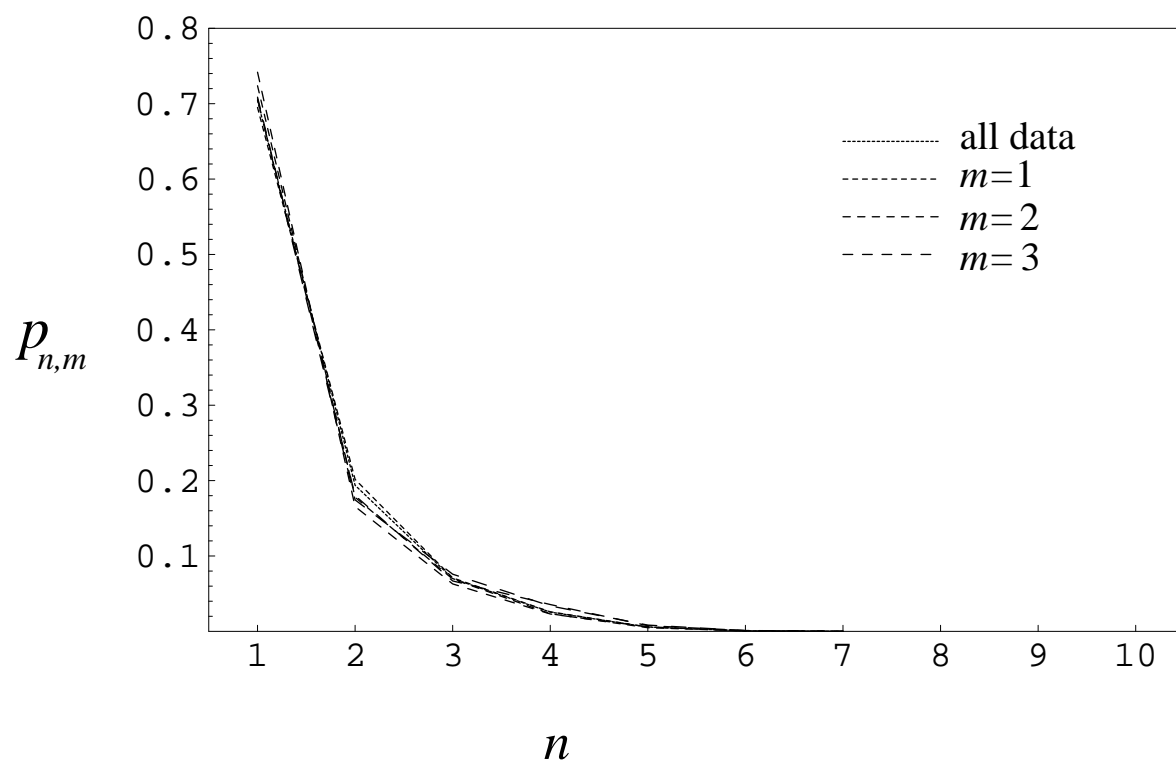


Figure 8: Plot of the probability distribution $p_{n,m}$ for George Eliot, *Middlemarch*

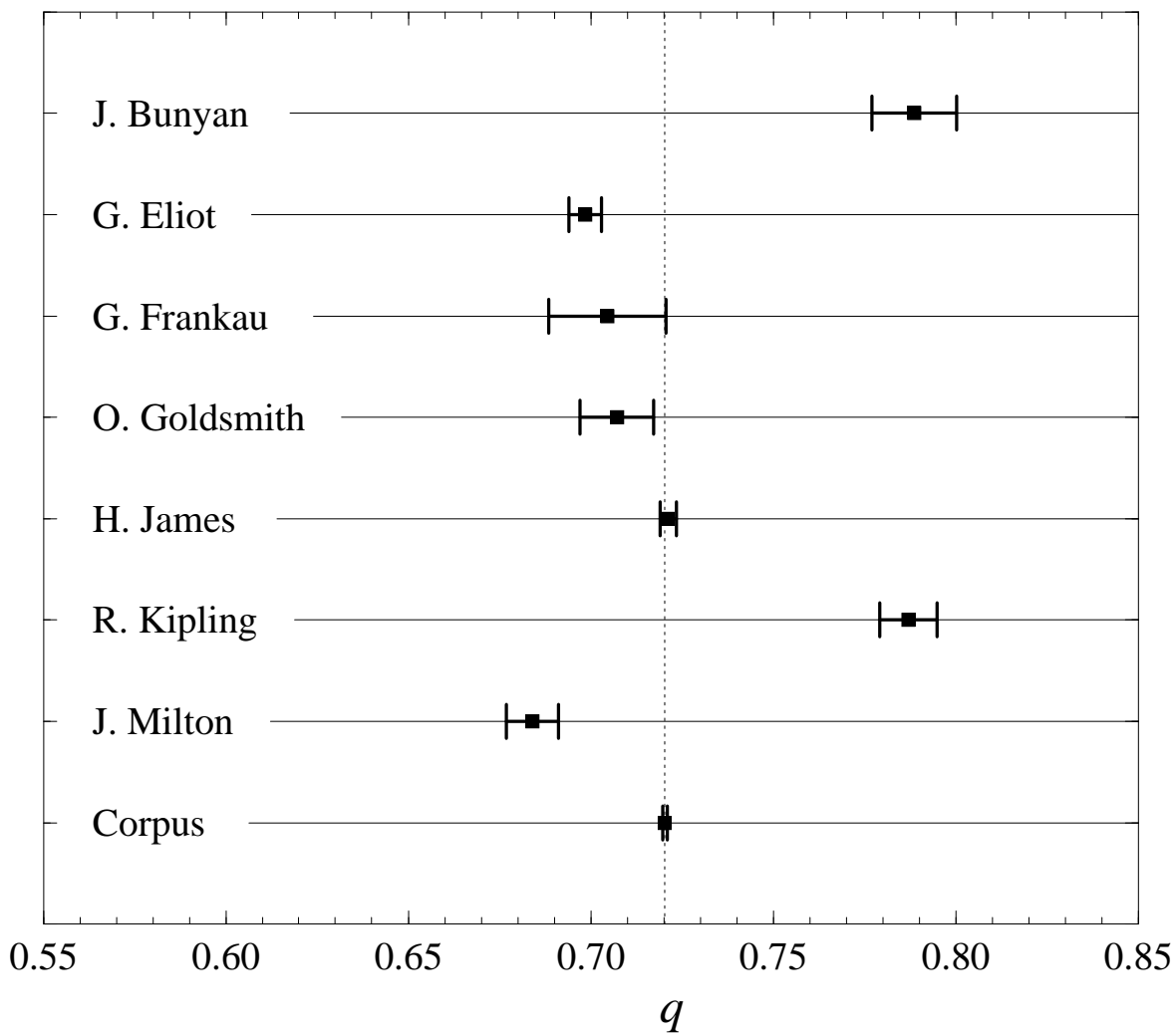


Figure 9: Plot of the value of q and its 3σ range for all the authors and the whole corpus

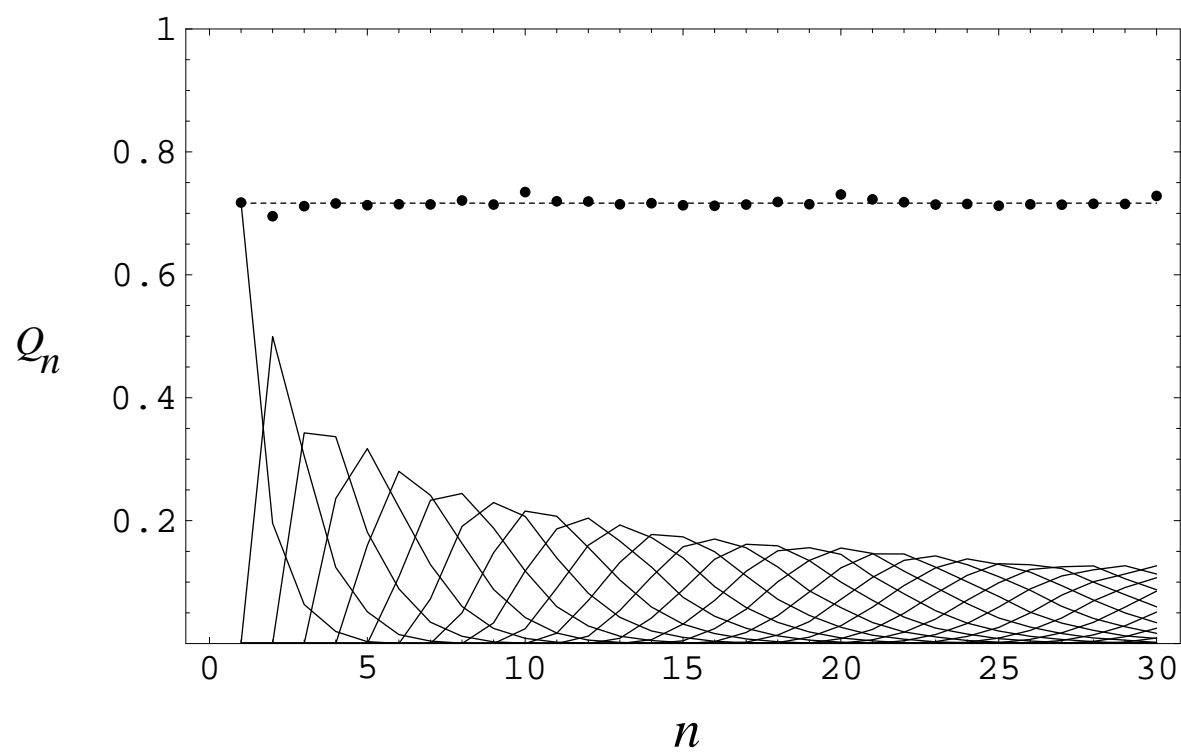


Figure 10: The frequency Q_n and the probabilities $P_{n,k}$ for Wordsworth's Prelude (57,570 words)

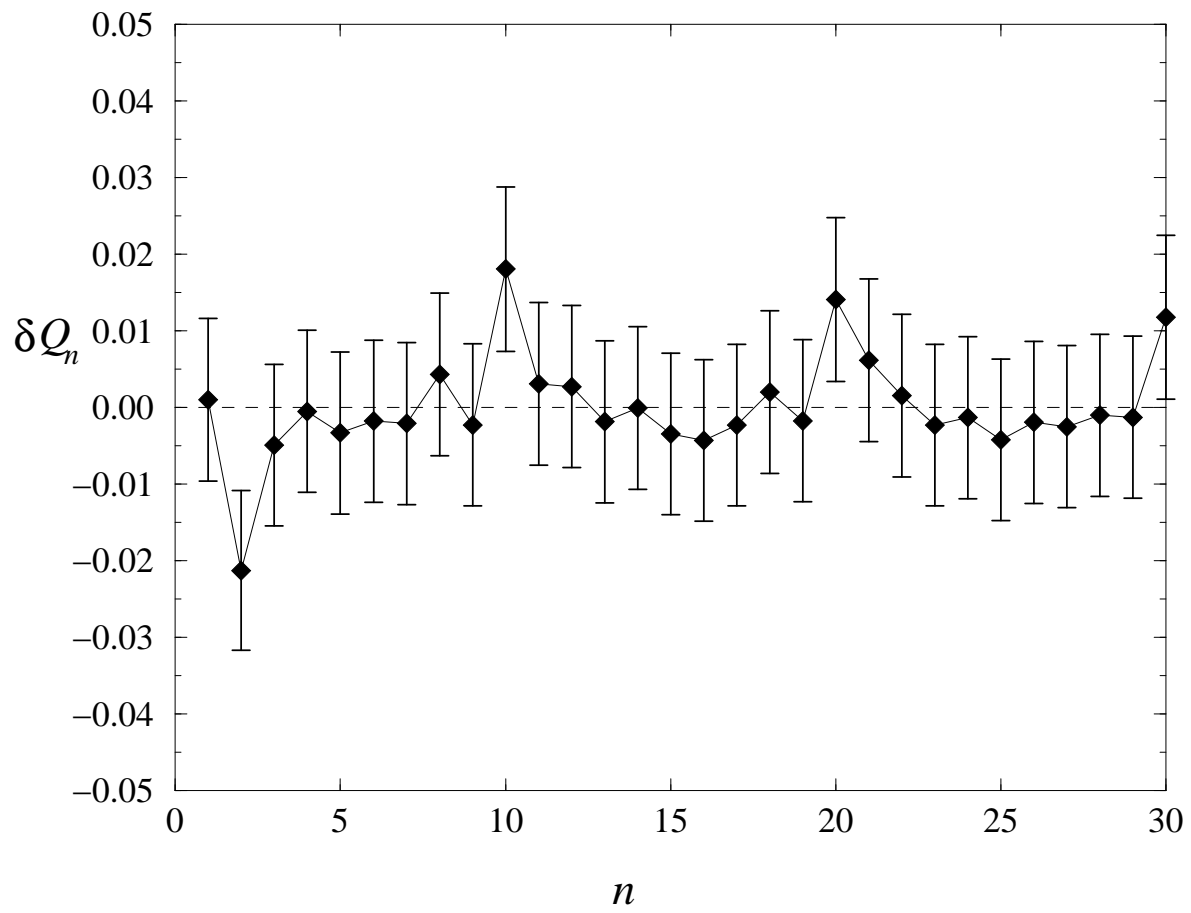


Figure 11: Detailed view of the frequency Q_n and the probabilities $P_{n,k}$ of Wordsworth's Prelude.